

Georgia State University  
**ScholarWorks @ Georgia State University**

---

Mathematics Theses

Department of Mathematics and Statistics

---

11-16-2006

# Changepoint Analysis of HIV Marker Responses

Joy Michelle Rogers

Follow this and additional works at: [https://scholarworks.gsu.edu/math\\_theses](https://scholarworks.gsu.edu/math_theses)

 Part of the [Mathematics Commons](#)

---

## Recommended Citation

Rogers, Joy Michelle, "Changepoint Analysis of HIV Marker Responses." Thesis, Georgia State University, 2006.  
[https://scholarworks.gsu.edu/math\\_theses/16](https://scholarworks.gsu.edu/math_theses/16)

This Thesis is brought to you for free and open access by the Department of Mathematics and Statistics at ScholarWorks @ Georgia State University. It has been accepted for inclusion in Mathematics Theses by an authorized administrator of ScholarWorks @ Georgia State University. For more information, please contact [scholarworks@gsu.edu](mailto:scholarworks@gsu.edu).

# CHANGEPOINT ANALYSIS OF HIV MARKER RESPONSES

by

JOY ROGERS

Under the Direction of Pulak Ghosh

## ABSTRACT

We will propose a random changepoint model for the analysis of longitudinal CD4 and CD8 T-cell counts, as well as viral RNA loads, for HIV infected subjects following highly active antiretroviral treatment. The data was taken from two studies, one of the Aids Clinical Group Trial 398 and one performed by the Terry Bein Community Programs for Clinical Research on AIDS. Models were created with the changepoint following both exponential and truncated normal distributions. The estimation of the changepoints was performed in a Bayesian analysis, with implementation in the WinBUGS software using Markov Chain Monte Carlo methods. For model selection, we used the deviance information criterion (DIC), a two term measure of model adequacy and complexity. DIC indicates that the data support a random changepoint model with the changepoint following an exponential distribution. Visual analyses of the posterior densities of the parameters also support these conclusions.

INDEX WORDS: CD4 count; CD8 count; Viral RNA load; DIC; Changepoint analysis

CHANGEPOINT ANALYSIS OF HIV MARKER RESPONSES

by

JOY ROGERS

A Thesis Submitted in Partial Fulfillment of the Requirements for the Degree of

Master of Science

in the College of Arts and Sciences

Georgia State University

2006

Copyright by  
Joy Michelle Rogers  
2006

# CHANGEPOINT ANALYSIS OF HIV MARKER RESPONSES

by

JOY MICHELLE ROGERS

Major Professor: Pulak Ghosh  
Committee: Jun Han  
Xu Zhang  
Yu-Sheng Hsu

Electronic Version Approved:

Office of Graduate Studies  
College of Arts and Sciences  
Georgia State University  
December 2006

## ACKNOWLEDGEMENTS

There are many people I owe deep gratitude to whose support assisted me in completing this thesis and my degree.

First of all, I would like to thank my advisor and committee chair, Dr. Pulak Ghosh, for his guidance, patience and suggestions. I would also like to thank the other members of my thesis committee:

I would also like to thank the other professors in the Department of Mathematics and Statistics for their instruction during my tenure at Georgia State University.

I also owe gratitude to my employer and coworkers for the flexibility and support offered in regards to my schedule.

Finally, I would like to thank my family for their never-ending patience and encouragement during the time I have spent in pursuit of this degree.

## TABLE OF CONTENTS

ACKNOWLEDGEMENTS.....	iv
LIST OF TABLES.....	vi
LIST OF FIGURES.....	vii
CHAPTER	
1 INTRODUCTION	1
2 CHANGEPOINT ANALYSIS.....	4
3 DESCRIPTION OF DATA SETS.....	6
4 DEVIANCE INFORMATION CRITERION FOR MODEL COMPARISON..	8
5 METHODS.....	10
6 RESULTS.....	12
7 DISCUSSION.....	24
REFERENCES.....	25
APPENDIX A.....	27

## LIST OF TABLES

Table 1. Model comparison with DIC values for data set one.....	13
Table 2. Model comparison with DIC values for data set two.....	13
Table 3. Parameter estimates, posterior means, and 95% posterior intervals for Models 1 and 2	18
Table 4. Parameter estimates, posterior means, and 95% posterior intervals for Model 3 and 4	18
Table 5. Parameter estimates, posterior means, and 95% posterior intervals for Models 5 and 6	18
Table 6. Parameter estimates, posterior means, and 95% posterior intervals for Models 7 and 8	18



## LIST OF FIGURES

Figure 1. Plot of the posterior mean values of the changepoints $K_i$ of model 1.....	14
Figure 2. Plot of the posterior mean values of the changepoints $K_i$ of model 2.....	14
Figure 3. Plot of the posterior mean values of the changepoints $K_i$ of model 3.....	15
Figure 4. Plot of the posterior mean values of the changepoints $K_i$ of model 4.....	15
Figure 5. Plot of the posterior mean values of the changepoints $K_i$ of model 5.....	16
Figure 6. Plot of the posterior mean values of the changepoints $K_i$ of model 6.....	16
Figure 7. Plot of the posterior mean values of the changepoints $K_i$ of model 7.....	17
Figure 8. Plot of the posterior mean values of the changepoints $K_i$ of model 8.....	17
Figures 9 – 11. Posterior density plots of $\mathbf{b}_1$ , $\mathbf{b}_2$ , and $\mathbf{b}_3$ for Model 1.....	19
Figures 12 – 14. Posterior density plots of $\mathbf{b}_1$ , $\mathbf{b}_2$ , and $\mathbf{b}_3$ for Model 2.....	20
Figures 15 – 17. Posterior density plots of $\mathbf{b}_1$ , $\mathbf{b}_2$ , and $\mathbf{b}_3$ for Model 3.....	20
Figures 18 – 20. Posterior density plots of $\mathbf{b}_1$ , $\mathbf{b}_2$ , and $\mathbf{b}_3$ for Model 4.....	21
Figures 21 – 23. Posterior density plots of $\mathbf{b}_1$ , $\mathbf{b}_2$ , and $\mathbf{b}_3$ for Model 5.....	21
Figures 24 – 26. Posterior density plots of $\mathbf{b}_1$ , $\mathbf{b}_2$ , and $\mathbf{b}_3$ for Model 6.....	22
Figures 27 – 29. Posterior density plots of $\mathbf{b}_1$ , $\mathbf{b}_2$ , and $\mathbf{b}_3$ for Model 7.....	22
Figures 30 – 32. Posterior density plots of $\mathbf{b}_1$ , $\mathbf{b}_2$ , and $\mathbf{b}_3$ for Model 8.....	23

## INTRODUCTION

Human Immunodeficiency Virus (HIV) is a retrovirus that uses the enzyme reverse transcriptase to transcribe RNA into DNA. Because transcription in the body generally occurs in the opposite direction – from DNA to RNA – retroviruses are highly susceptible to mutations. As a result, these viruses are very difficult to treat with antibiotics and thus far impossible to develop vaccines for (<http://en.wikipedia.org/wiki/Hiv>). Thus, it is important to study HIV's infection mechanisms in order to develop drug therapy programs that can slow the destruction of the body.

HIV assaults the immune system by attacking the body's cells, primarily CD4+ Helper T cells. The term “Helper” is used because of the function of the CD4 cell; once the cells are activated, they divide quickly and emit cytokines, which are proteins that help the immune response. Along with CD4+ T cells, levels of CD8 T cells can also be used as a marker for HIV progression. CD8 T cells, unlike CD4 cells, are not “Helper” T cells. Instead, they are Cytotoxic T cells, which means they destroy cells that are virally infected ([http://en.wikipedia.org/wiki/T\\_cells](http://en.wikipedia.org/wiki/T_cells)). The final key definition is that of viral RNA. As mentioned above, HIV is a retrovirus and contains RNA as the hereditary material. This RNA contained in the retroviral particle is called viral RNA, since it is the template for further reproduction of viral cells (Kimball, 2005).

Because the CD4+ cells are the primary target of HIV, the level of these cells in the body is often used as an indicator of the advancement of HIV infection (Wang, 2006). The levels of these cells in conjunction with the viral load (number of copies of HIV-1 RNA per millimeter of plasma) are widely used to predict progression of the disease to full-blown Acquired Immune Deficiency Syndrome (AIDS) and eventually, death (Kiuchi, 1995). Since the

levels of these markers change from day-to-day, it is advisable to look at the general pattern over time in order to model the advancement of the disease.

As the virus progresses, infected CD4+ cells are destroyed. However, the body is simultaneously producing new CD4+ cells in order to produce antibodies and attempt to attack the virus. It is generally accepted that immediately after infection, CD4+ cells rapidly multiply as an initial line of defense against HIV; however, if the disease is left untreated, this increase in T-cell levels soon begins to decline as the virus attacks them (Gumel, 2001). This decline can last up to several years, and often steepens just before the patient is diagnosed with AIDS.

This progression can be slowed or reversed when the patient is treated. When an individual begins on a drug program such as Highly Active Anti-Retroviral Therapy (HAART), the immune system can be restored to some extent because the drugs suppress the replication of the HIV cells and reduce the viral load (Ghosh and Vaida, 2006). As a result of the medication, the CD4+ counts will begin to rise in a two stage process. The first stage is a steep increase in CD4+ counts; the second, a more gradual increase (Deeks, 2004). At least one historical study of this problem has estimated the time of the first increase to be around 12 weeks (Hunt, et al., 2003); however, that model did not allow for variability between subjects, which can be done using a Bayesian approach to changepoint analysis. This paper is concerned with estimating the changepoint for each patient at the time where the steep increase slows to the gradual increase of CD4+ levels. We also extend our method to two other markers of HIV/AIDS – CD8 T-cells and viral RNA load.

This thesis is organized in the following manner: Chapter 2 discusses changepoint analysis and its background as a statistical process. This is followed by a description of the data sets used in Chapter 3, and a description in Chapter 4 of the measure used in model selection, the

Deviance Information Criterion. Chapter 5 explains the methods used in testing different models, and Chapter 6 details the results of the testing. The thesis concludes with a discussion of these results.

## CHANGEPOINT ANALYSIS

In statistical analysis, it is often important to be able to determine if and when a change occurred in the behavior of data, especially time ordered data. For example, it may be important to analyze a series of counts of violent crimes before and after a specific crime-fighting program was introduced to determine if the crime rate decreased (Loschi, et al, 2005). Historically, methods such as control charts have been used to detect changes in data; however, these charts often fall short in detecting change if the change in the data is subtle. Also, control charts control the point-wise error rate, as opposed to the change-wise error rate, so if a data set is large, the control chart may indicate that a point is outside the control limits even when no change has occurred, or vice-versa.

Changepoint analysis is a powerful alternative method used to detect even slight changes in the distribution of time-ordered data. Changepoint analysis is especially beneficial because it can provide confidence levels and confidence intervals for the estimate of the point where the change occurs. This helps to ensure that each change detected is real. Also, changepoint analysis is reasonably robust to outliers, and can be adjusted to different types of data sets (attribute data, individual values, counts, and standard deviations, for example) easily. Finally, because changepoint analysis automates the process of calculating new control limits following each change, it is much easier to use when dealing with large sets of historical data (Taylor, 2000).

There are some shortcomings when considering changepoint analysis. Unlike some other tests, including control charts, it does not detect isolated abnormal points. Also, because of the way changepoint analysis is performed (using a bootstrapping approach), it will not produce identical results each time it is performed due to the random selection of bootstrap samples. This

can be overcome by increasing the number of iterations. However, these inadequacies do not affect the use of the analysis in this paper, since approximately 50,000 iterations were used for each model.

Changepoint analysis has been an active area of research to model progression of HIV and AIDS. It is often used in determining the behavior of various biological markers, like CD4 and CD8 T-cell counts and viral RNA load surrounding major events in the advancement of the disease. Multiple changepoint testing has been used by Halpern to assess the possibility of genetic recombination in HIV sequences (2000). Kiuchi, et al (1995) used a Bayesian changepoint analysis (specifically using an EM algorithm and MCMC) to estimate the time that a rapid decline of CD4 cells begins just prior to diagnosis with AIDS at approximately 1 year (with a standard deviation of 9 months). Similarly, Lange, et al. (1992) developed models to find posterior distributions of time for CD4 T-cell numbers to reach a specified level (for example, near seroconversion or progression to AIDS).

For our purposes, we use a Bayesian approach to create a model that allows variability between patients in their CD4 levels and changepoint times. Each patient's profile will be different (including CD4, CD8, and RNA counts), and the models used here will allow for random changepoints (Ghosh and Vaidya, 2006).

## DESCRIPTION OF DATA SETS

Two data sets are used in the construction of the model – AIDS Clinical Trial Group (ACTG) 398 and a set taken from a study conducted by the Terry Bein Community Programs for Clinical Research on AIDS.

ACTG 398 was a study completed on 481 subjects with advanced stages of HIV infection. The study was randomized, double-blind, and placebo-controlled, and compared four different highly active antiretroviral treatment (HAART) regimens. All regimens in the study contained these antiretroviral drugs: abacavir (a nucleoside reverse transcriptase inhibitor – NRTI), adefovir dipivoxil (NRTI), efavirenz (a non-NRTI – NNRTI), and amprenavir (a protease inhibitor –PI). Three of the four study groups were given a second PI (either saquinavir, indinavir, or nelfinavir), and the fourth was given a placebo. The study was designed to compare the proportion of patients who had “virologic failure after 24 weeks of study between the double-PI arms and the single-PI arm” (Ghosh and Vaida, 2006). The data set used in this study contains measurements of CD4 cells, CD8 cells, and viral RNA loads for the subjects at weeks 0, 2, 4, 8, 16, 24, 32, 40, and 48.

The second data set, as mentioned previously, was taken from a study conducted by the Terry Bein Community Programs for Clinical Research on AIDS. The study was a multicenter, randomized, open-label, community based clinical trial of 467 patients. These criteria required for inclusion into the study were as follows: at least 13 years of age; “having an AIDS defining condition or two CD4 lymphocyte counts of 300 cells or less per cubic millimeter, with either a positive serologic test for HIV or a working diagnosis of HIV infection” (Abrams, et al, 1994); and they had previously attempted (unsuccessfully) therapy with the drug zidovudine. This drug must have led to either intolerance of the drug or progression of disease during therapy. The trial

was designed to test the performance of two other drugs, didanosine or zalcitabine, and to see which of these dideoxynucleoside monotherapies performed better after zidovudine failed. The patients were randomly assigned to receive either didanosine or zalcitabine, and allowed to switch from one drug to the other once during the course of the trial if intolerance developed. 230 patients were assigned to didanosine, and 237 to zalcitabine, and the results of the study concluded that zalcitabine was at least as efficacious as didanosine in delaying progression and death, although neither drug offered substantial long-term benefit. Further results and analysis are reported in the article by Abrams, et al. The specific data set used here was taken from Guo & Carlin (2004), using the CD4 counts recorded at study entry, and again at the 2, 6, 12, and 18 month visits.



## DEVIANE INFORMATION CRITERION FOR MODEL SELECTION

It is necessary to have some criteria to select the best model when multiple models are tested. In this study, we use deviance information criterion (DIC), which was proposed in 2002 by Spiegelhalter, Best, Carlin, and van der Linde, and discussed by Berg, et al (2004). This is a generalized version of the Akaike information criterion (AIC) and similar to the Bayesian information criterion (BIC), and like these, it compares a model based on measures of its adequacy and its complexity. In order to accomplish this, the DIC uses two components: a Bayesian measure of model fit,  $\overline{D}$ , and the effective number of parameters,  $p_D$ . The goodness of fit term,  $\overline{D}$ , is defined as the posterior expectation of the deviance:

$$\overline{D} = E_{\mathbf{q}|\mathbf{y}}[D(\mathbf{q})] = E_{\mathbf{q}|\mathbf{y}}[-2 \ln f(\mathbf{y} | \mathbf{q})].$$

$\overline{D}$  is smaller for better-fitting models.

The second term,  $p_D$ , is defined as the difference between the posterior mean of the deviance and the deviance evaluated at the posterior mean  $\overline{\mathbf{q}}$  of the parameters:

$$p_D = \overline{D} - D(\overline{\mathbf{q}}) = E_{\mathbf{q}|\mathbf{y}}[D(\mathbf{q})] - D(E_{\mathbf{q}|\mathbf{y}}[\mathbf{q}]) = E_{\mathbf{q}|\mathbf{y}}[-2 \ln f(\mathbf{y} | \mathbf{q})] + 2 \ln f(\mathbf{y} | \overline{\mathbf{q}}).$$

This means that  $p_D$  can be understood as the expected excess of the true over the estimated residual information in data  $\mathbf{y}$  conditional on  $\mathbf{q}$ ; in other words,  $p_D$  is the expected reduction in uncertainty due to estimation.

Putting the above two terms together and rearranging, we can define DIC as follows:

$$\text{DIC} = D(\overline{\mathbf{q}}) + 2p_D.$$

Myung, et al. (2005) gives three reasons why DIC is a desirable criterion for model selection in the Bayesian context. The first is that, as mentioned above, the DIC is a Bayesian counterpart of the AIC (in fact, the two are asymptotically equivalent for normal linear models).

DIC improves over AIC in that its penalty term measures the complexity of order-constrained models, whereas AIC cannot discriminate between models with the same number of parameters but differing implementation of order-constraints. Secondly, DIC does not require nesting of competing models. Also, models do not have to be assumed true. The third reason DIC works well is that it can be estimated by the posterior parameter samples generated by any Monte Carlo algorithm. This allows DIC to be a superior criterion to the Bayes factor, another model selection standard (it is defined as the ratio of the marginal likelihood under one model to the marginal likelihood under a competing model). The Bayes factor is very sensitive to the assumed prior distribution and can be difficult to compute and interpret for hierarchical models. By comparison, computing DIC via MCMC is almost trivial (Ghosh and Vaidya, 2006).

## METHODS

This paper will consider a model of the form

$$y_{ij} = \mathbf{b}_1 + \mathbf{b}_2(t_{ij} - K_i)_- + \mathbf{b}_3(t_{ij} - K_i)_+ + b_i + e_{ij}, \quad (1)$$

where  $y_i = (y_{i1}, \dots, y_{in_i})^T$  is the response vector for individual  $i$ , and  $t_i = (t_{i1}, \dots, t_{in_i})^T$  is the corresponding vector of observation times,  $i = 1, \dots, n$ . Also, let  $K_i$  be the changepoint for individual  $i$  and  $\mathbf{b} = (\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3)^T$  be the vector of fixed effects. Further define  $\mathbf{b}_1$  as the overall intercept, and let  $\mathbf{b}_2$  and  $\mathbf{b}_3$  represent the overall slope before and after the changepoint  $K_i$ . Let  $b_i$  be the random effects for individual  $i$ , and  $e_{ij}$  be the error. The random effects  $b_i$  are assumed to follow a normal distribution with mean 0 and standard deviation  $\mathbf{s}_b^2$ ; likewise, we assume that  $e_{ij}$  follows a normal distribution. These distributions can be written as

$$b_i \sim N(0, \mathbf{s}_b^2), e_{ij} \sim N(0, \mathbf{s}^2). \quad (2)$$

One can take a more general model than (1) that includes additional random effects, but it does not give much benefit for the data sets we are analyzing.

Conditionally on the changepoints  $K_i$ , (1) follows a mixed-effects model, which can be written in general as

$$y_i = X_i \mathbf{b} + Z_i b_i + e_i, \quad (3)$$

where  $X_i$  and  $Z_i$  are  $n_i \times p$  and  $n_i \times q$  matrices of covariates for subject  $i$ , corresponding to the fixed and random effects, respectively.  $X_i$  includes baseline covariates, possibly time-varying covariates, such as treatment adherence, and the before-changepoint and after-changepoint vectors  $(t - K_i)_-$  and  $(t - K_i)_+$ ; in general,  $Z_i$  is a submatrix of  $X_i$ . The random effects are given by (2), and the errors are assumed normal,  $e_i \sim N(0, \mathbf{s}^2 I_{n_i})$ , independently of each other.

The changepoints  $K_i$  are assumed to be random,

$$K_i \sim F(\cdot; \mathbf{g}), \quad (4)$$

independently of each other;  $F(\cdot; \mathbf{g})$  is a parametric family of distributions (in this case, specifically exponential and truncated normal). The equations (3), (4) define a general class of mixed effects random changepoint models for the biological markers.

It is also important to note that each of the covariates used was transformed in order to stabilize the variance, improve normality of errors, and linearity of the mean. The CD4 and CD8 cell counts were placed on a square-root scale, and the RNA counts were adjusted using a logarithmic transformation.

We fit the model using a fully Bayesian approach. The full model specification includes prior distributions for the parameter  $\mathbf{q}$ ,

$$f(y, T, \mathbf{q}) = f(y, T|\mathbf{q})f(\mathbf{q}). \quad (5)$$

We assume the unknown parameters  $\mathbf{b}$ ,  $\mathbf{s}_b^2$ ,  $\mathbf{s}^2$ , and  $\mathbf{g}$  to be mutually independent in the prior.

The hyperparameters in the prior distributions were chosen so that the priors are uninformative. We took

$$\mathbf{b}_j \sim N(0, 1000), \quad \text{independently} \quad (6)$$

$$\mathbf{s}^2 \sim IG(0.001, 0.001) \quad (7)$$

$$\mathbf{s}_b^2 \sim IG(0.001, 0.001) \quad (8)$$

$$\mathbf{g} \sim G(0.1, 0.1); \quad (9)$$

$x \sim IG(a, b)$  means that  $1/x$  has the Gamma distribution with mean  $a/b$  and variance  $a/b^2$ , and

$G(a, b)$  denotes a Gamma distribution with parameter  $a, b$ .

## RESULTS

For the first data set (from ACTG 398), we considered several models for  $y_{ij}$ , as follows:

- **Model 1:** Uses the CD4 cell level as the response, with the changepoint following an exponential distribution,  $K_i \sim \text{Exp}(\mathbf{I})$
- **Model 2:** Same as **Model 1**, but with the changepoint following a truncated normal distribution,  $K_i \sim \text{TN}(\mathbf{2}, \mathbf{0.1})\mathbf{I}(0,)$
- **Model 3:** Same as **Model 1**, but with CD8 level as the response
- **Model 4:** Same as **Model 2**, but with CD8 level as the response
- **Model 5:** Same as **Model 1**, but with  $\log_{10}$  HIV-1 RNA level as the response
- **Model 6:** Same as **Model 2**, but with  $\log_{10}$  HIV-1 RNA level as the response

For data set two, the same models as **Model 1** and **Model 2** were used. All models have a random changepoint for the responses.

We implemented the model in WinBUGS, using Markov chain Monte Carlo sampling from the posterior distribution (see Appendix A for further information on WinBUGS implementation). It is important to determine the number of iterations necessary to achieve convergence to the stationary distribution. For this study, it was concluded that 25,000 iterations were sufficient for the burn-in period.

Table 1 contains the model comparison for the different models used, including the DIC value, the effective number of parameters (or effective degrees of freedom)  $p_D$ , and an estimate for  $\bar{K}$ , the posterior mean of the changepoint, and its standard deviation for the 455 subjects, in weeks. Models 1 – 6 deal with the first data set (from ACTG 398). Models 7 and 8 use the data from set two. For data set 1, the two models using the CD8 counts (models 3 and 4) have the highest DIC values at 15689.4 and 15791.8, respectively, so they have the least support. Using

CD4 as the response reduces the DIC values to 12328.9 for model 1 and 12332.3 for model 2.

There is a large improvement, however, when using RNA as the response, and we get 8075.3 and 8431.5 as the DIC values for models 5 and 6 respectively. The number of effective degrees of freedom  $p_D$  is smallest for model 3. All of the models confirm that when  $K_i$  follows the exponential distribution, the model is more adequate. All six models for data set 1 have similar posterior means of the changepoint values, which is placed, on average, at 3.7 weeks, with standard deviation of 0.39 weeks.

	DIC	$p_D$	$\bar{K}$ (SD)
Model 1	12328.9	665.5	3.769 (0.03)
Model 2	12332.3	692.5	3.405 (0.11)
Model 3	15689.4	551.3	3.964 (0.18)
Model 4	15791.8	562.6	3.408 (0.11)
Model 5	8075.3	588.4	4.364 (0.21)
Model 6	8431.5	612.4	3.448 (0.11)

**Table 1. Model comparison using DIC values for data set 1.**

Similar to data set 1, for the second data set, model 7 (which uses the exponential distribution for the changepoints) is superior to model 8 (which uses the truncated normal distribution). Model 7 does have a smaller  $p_D$  value at 375.5. The posterior mean of the changepoint value is placed, on average, at 3.7 weeks, with a standard deviation of 0.43 weeks.

	DIC	$p_D$	$\bar{K}$ (SD)
Model 7	6117.6	375.5	3.988 (0.18)
Model 8	6507.5	509.2	3.375 (0.10)

**Table 2. Model Comparison using DIC values for data set 2.**

For all models, DIC gives credible values with positive and meaningful  $p_D$  values.

Figures 1 through 8 below show the density plots of the posterior means for the changepoints  $K_i$  for each of the models.

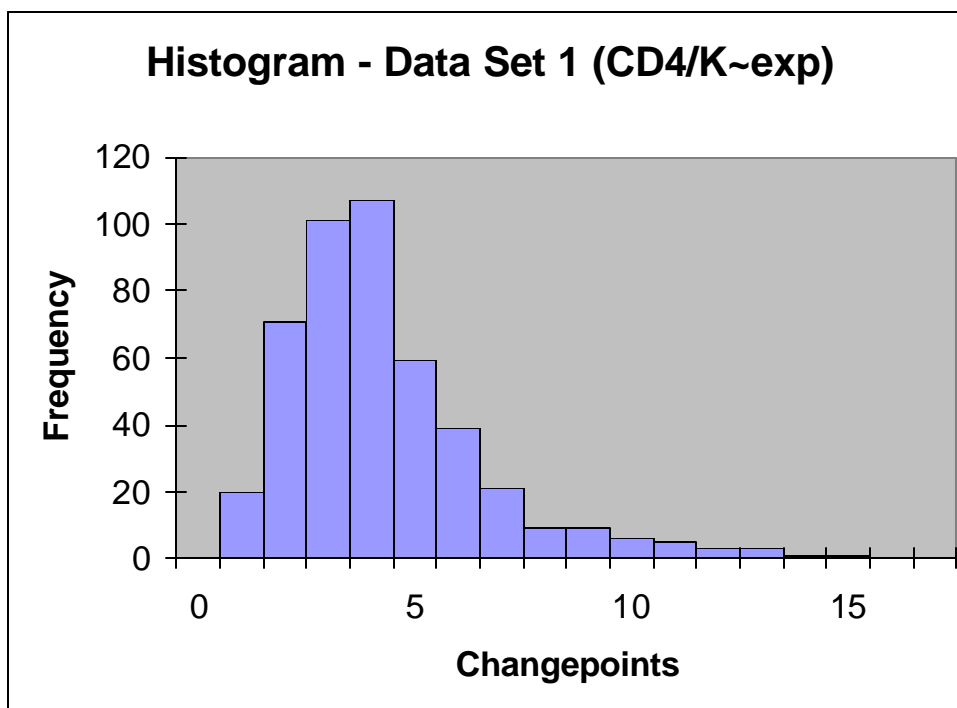


Figure 1. Plot of the posterior mean values of the changepoints  $K_i$  of model 1

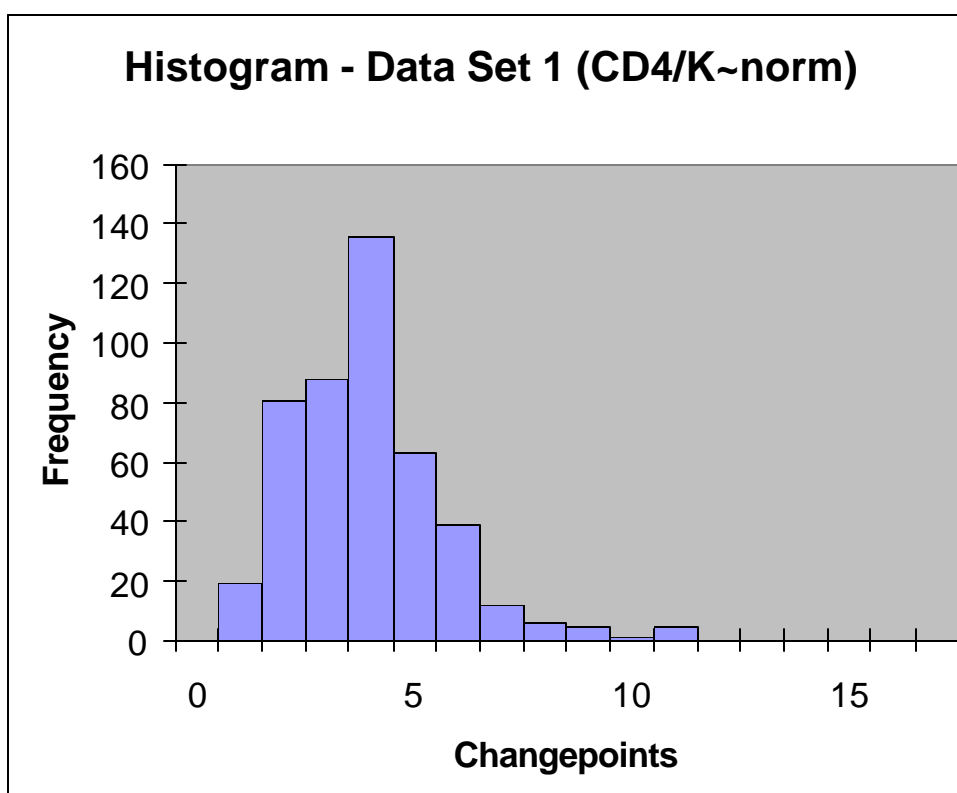


Figure 2. Plot of the posterior mean values of the changepoints  $K_i$  of model 2

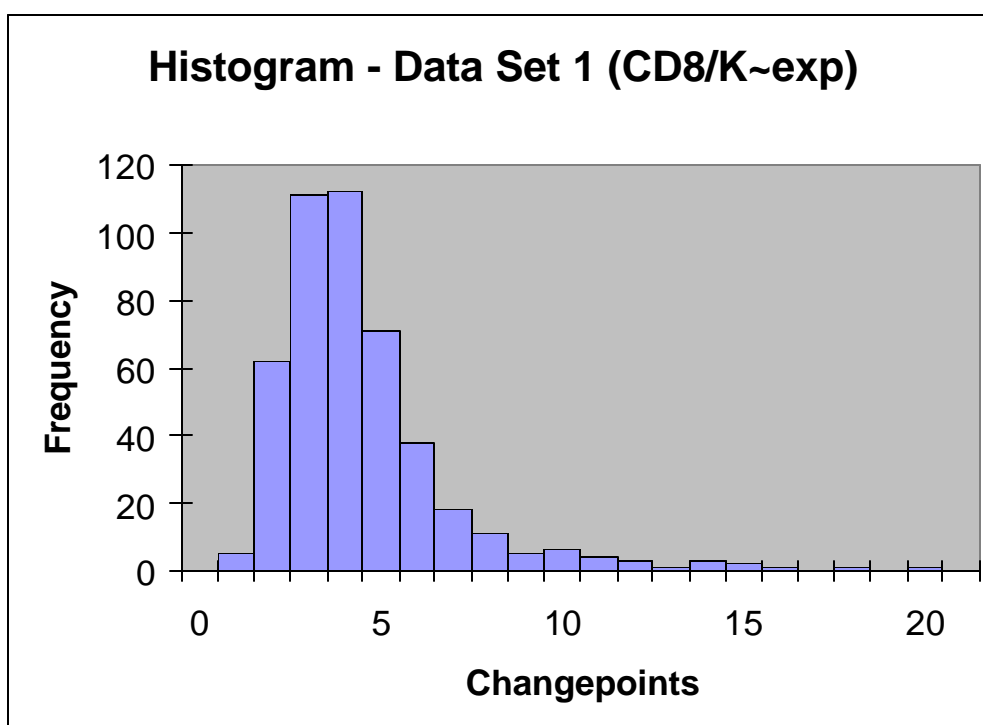


Figure 3. Plot of the posterior mean values of the changepoints  $K_i$  of model 3

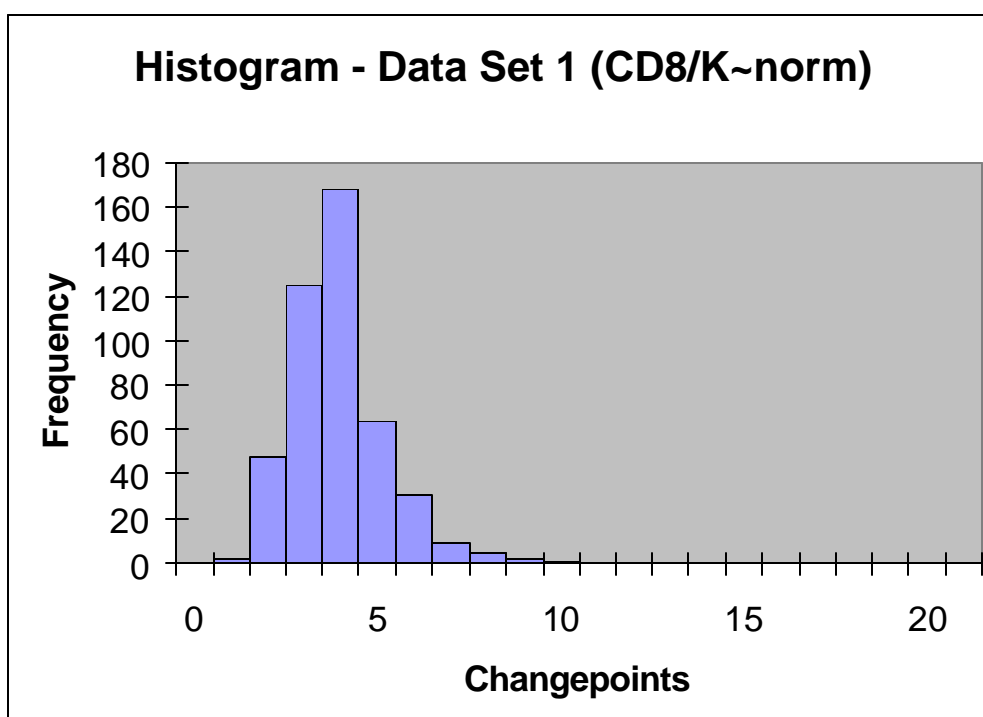


Figure 4. Plot of the posterior mean values of the changepoints  $K_i$  of model 4



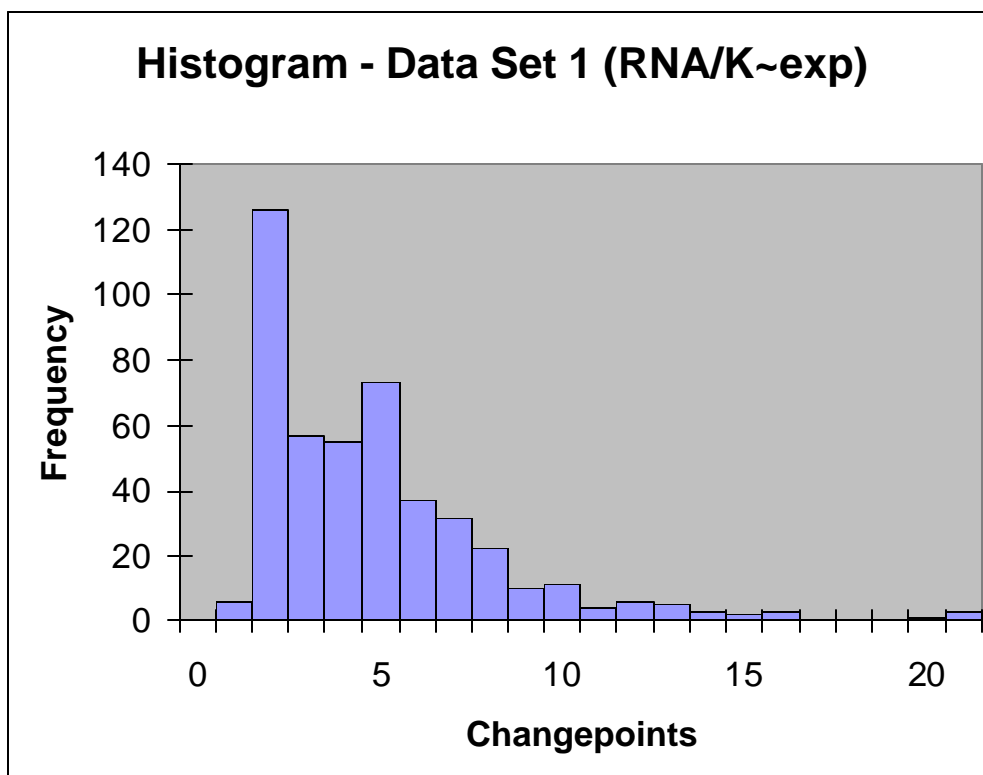


Figure 5. Plot of the posterior mean values of the changepoints  $K_i$  of model 5

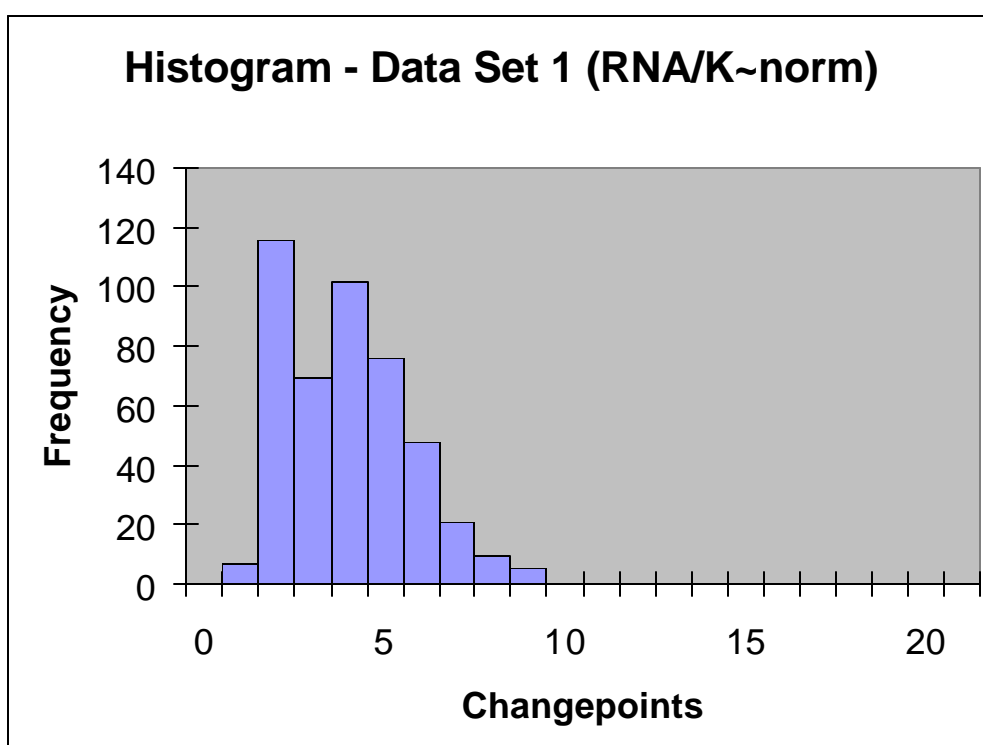


Figure 6. Plot of the posterior mean values of the changepoints  $K_i$  of model 6

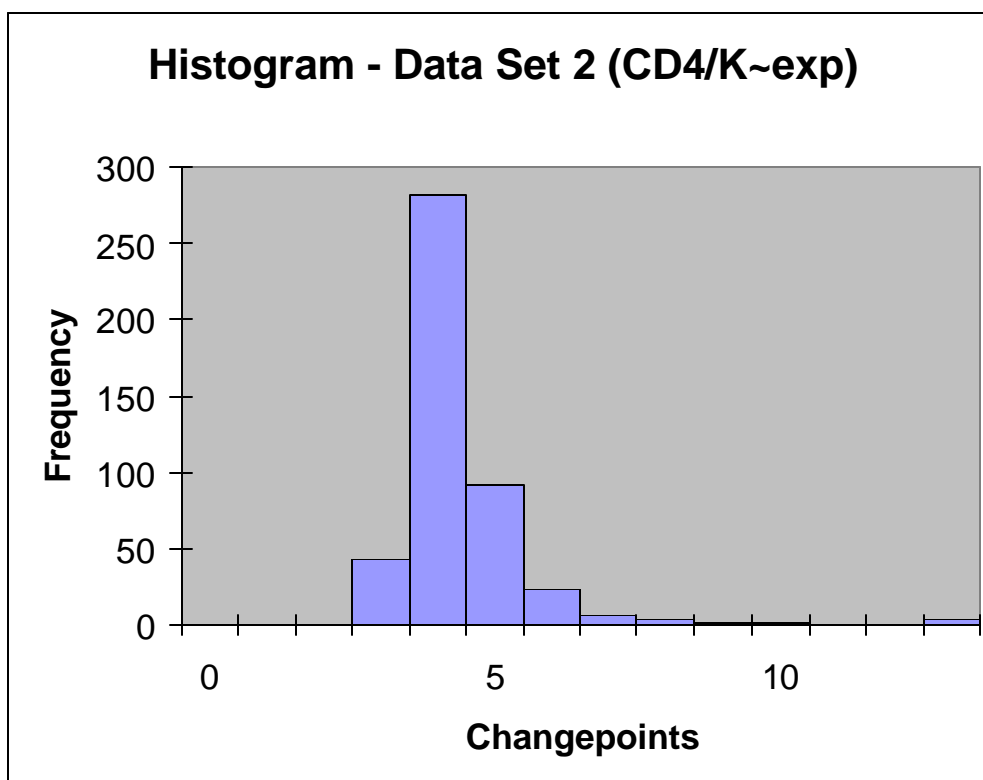


Figure 7. Plot of the posterior mean values of the changepoints  $K_i$  of model 7

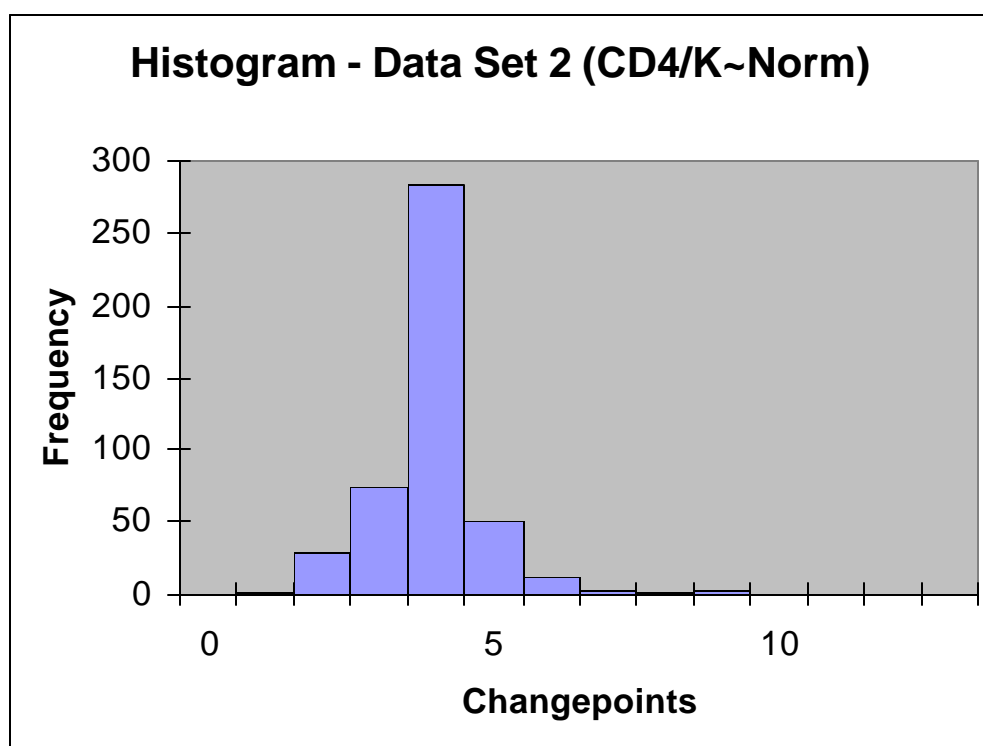


Figure 8. Plot of the posterior mean values of the changepoints  $K_i$  of model 8

The parameter estimates for each model are shown in Tables 2-5.

Model 1			Model 2		
Parameter Estimate	Posterior Mean	95% Posterior Interval	Parameter Estimate	Posterior Mean	95% Posterior Interval
$\beta_1$	15.19	(14.7, 15.67)	$\beta_1$	14.97	(14.47, 15.45)
$\beta_2$	0.021	(0.017, 0.025)	$\beta_2$	-0.03345	(-0.037, -0.030)
$\beta_3$	-0.070	(-0.082, -0.058)	$\beta_3$	0.1149	(0.102, 0.129)
$\sigma_b$	5.168	(4.832, 5.528)	$\sigma_b$	5.129	(4.795, 5.488)
$\sigma$	1.959	(1.897, 2.023)	$\sigma$	1.855	(1.798, 1.915)

Table 3. Parameter estimates, posterior means, and 95% posterior intervals for Models 1 and 2

Model 3			Model 4		
Parameter Estimate	Posterior Mean	95% Posterior Interval	Parameter Estimate	Posterior Mean	95% Posterior Interval
$\beta_1$	29.4	(28.74, 30.04)	$\beta_1$	29.73	(29.09, 30.38)
$\beta_2$	-0.029	(-0.034, -0.023)	$\beta_2$	0.032	(0.025, 0.039)
$\beta_3$	0.047	(0.028, 0.066)	$\beta_3$	-0.162	(-0.187, -0.136)
$\sigma_b$	6.711	(6.263, 7.195)	$\sigma_b$	6.701	(6.252, 7.186)
$\sigma$	3.658	(3.543, 3.775)	$\sigma$	3.719	(3.603, 3.838)

Table 4. Parameter estimates, posterior means, and 95% posterior intervals for Models 3 and 4

Model 5			Model 6		
Parameter Estimate	Posterior Mean	95% Posterior Interval	Parameter Estimate	Posterior Mean	95% Posterior Interval
$\beta_1$	3.977	(3.83, 4.07)	$\beta_1$	3.952	(3.86, 4.046)
$\beta_2$	0.006	(0.005, 0.006)	$\beta_2$	0.009	(0.008, 0.010)
$\beta_3$	-0.027	(-0.030, -0.024)	$\beta_3$	-0.035	(-0.038, -0.031)
$\sigma_b$	0.928	(0.861, 1.001)	$\sigma_b$	0.912	(0.844, 0.984)
$\sigma$	0.765	(0.744, 0.787)	$\sigma$	0.772	(0.751, 0.793)

Table 5. Parameter estimates, posterior means, and 95% posterior intervals for Models 5 and 6

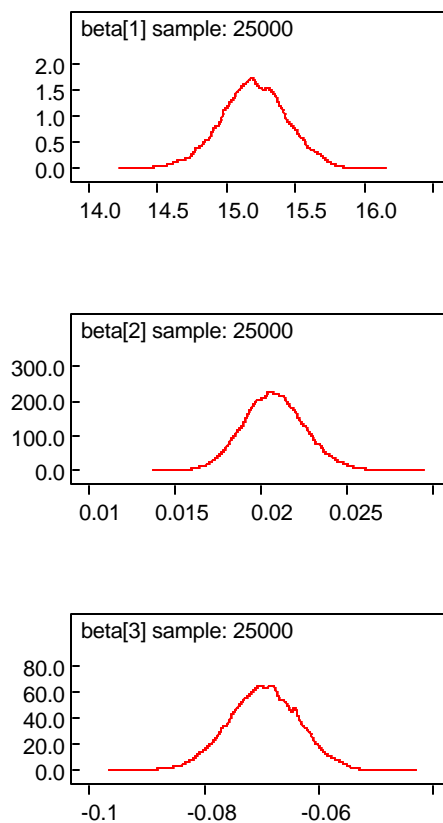
Model 7			Model 8		
Parameter Estimate	Posterior Mean	95% Posterior Interval	Parameter Estimate	Posterior Mean	95% Posterior Interval
$\beta_1$	7.088	(6.598, 7.587)	$\beta_1$	7.171	(6.756, 7.607)
$\beta_2$	0.004	(-0.046, 0.056)	$\beta_2$	0.059	(0.036, 0.081)
$\beta_3$	-0.220	(-0.352, -0.086)	$\beta_3$	-0.373	(-0.446, -0.300)
$\sigma_b$	4.512	(4.208, 4.833)	$\sigma_b$	4.537	(4.239, 4.862)
$\sigma$	1.86	(1.766, 1.957)	$\sigma$	1.841	(1.746, 1.942)

Table 6. Parameter estimates, posterior means, and 95% posterior intervals for Models 7 and 8 (data set 2)

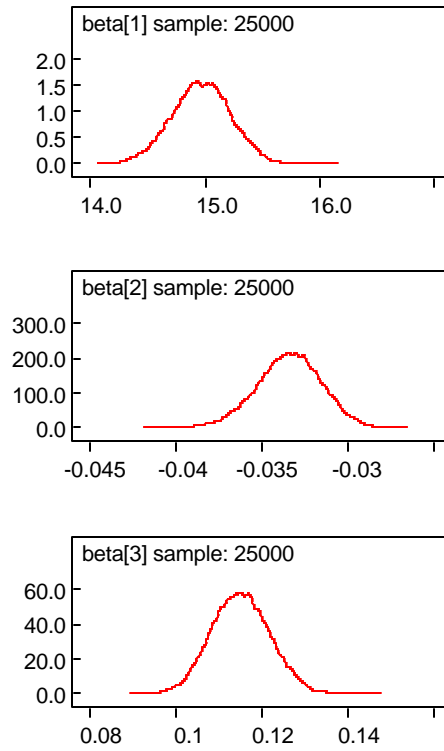
Each pair of models gives similar estimates for the parameters. Of particular interest is Table 4, which gives the estimates for Models 5 and 6, the models for data set 1 that gave the lowest DIC values. Specifically, in Model 5, the initial slope ( $b_2$ ) is of about 0.006 (in  $\log_{10}$

RNA per week), and is significantly positive, since the 95% posterior interval does not contain 0. The subsequent slope ( $b_3$ ), following the changepoint, is much smaller, as expected. Its negative value indicates that the longitudinal trajectory is headed down after the changepoint. The same observations hold true for model 6. For the second data set, we see in Table 5 that similar conclusions can be drawn about models 7 and 8.

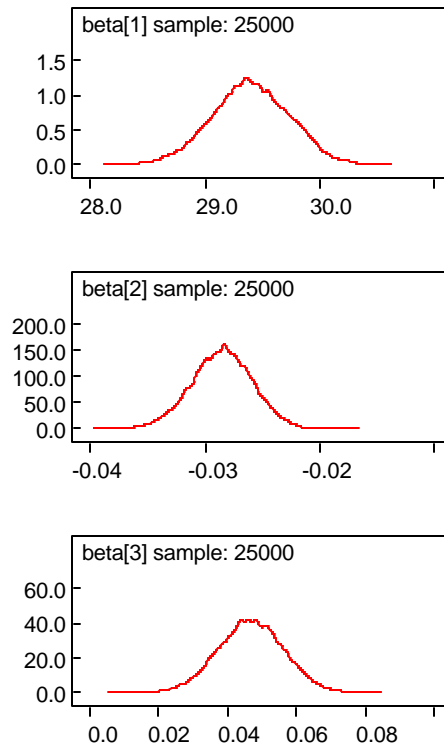
We can also perform a visual inspection the beta parameters by looking at the posterior density plots, seen in the figures below.



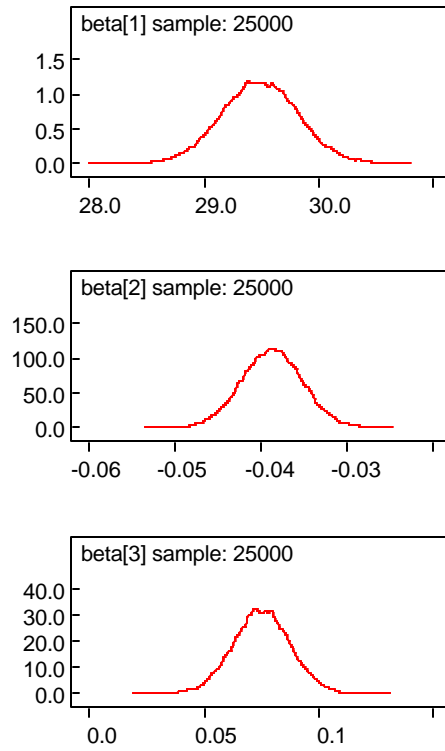
**Figures 9 - 11: Posterior density plots of  $b_1$ ,  $b_2$ , and  $b_3$  for Model 1**



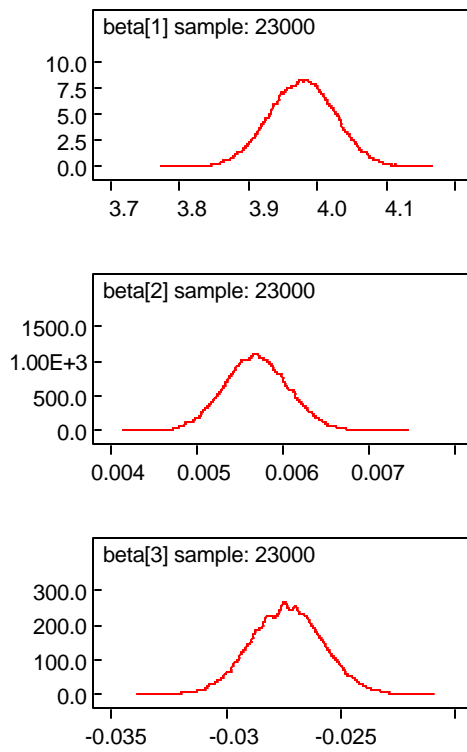
**Figures 12-14: Posterior density plots of  $b_1$ ,  $b_2$ , and  $b_3$  for Model 2**



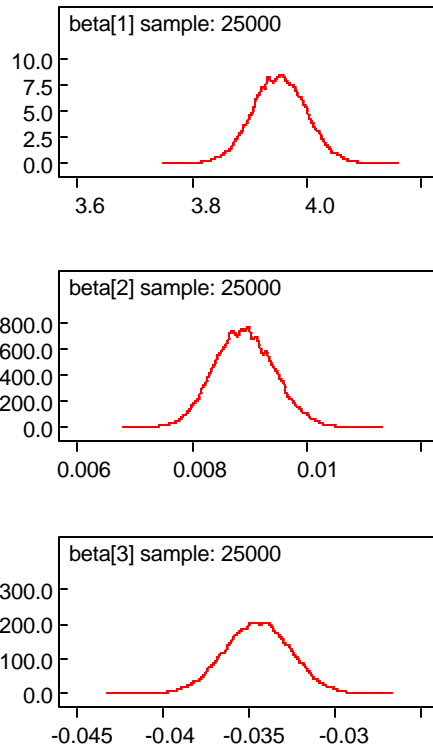
**Figures 15-17: Posterior density plots of  $b_1$ ,  $b_2$ , and  $b_3$  for Model 3**



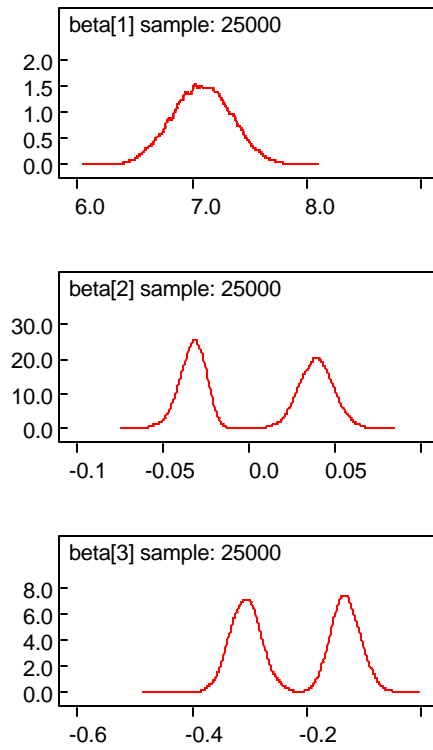
**Figures 18-20: Posterior density plots of  $b_1$ ,  $b_2$ , and  $b_3$  for Model 4**



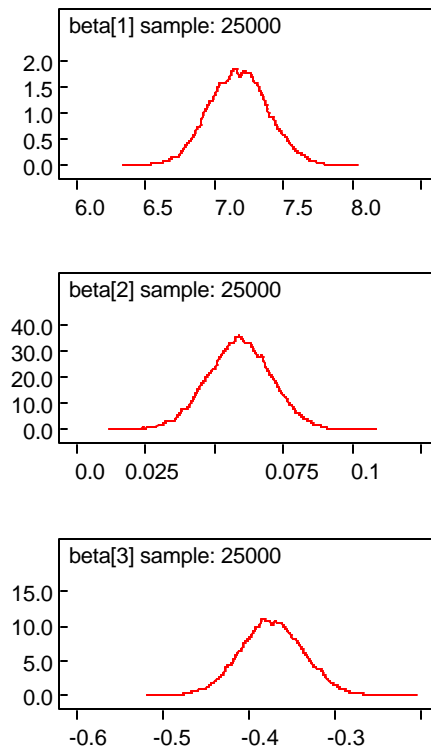
**Figures 21 - 23: Posterior density plots of  $b_1$ ,  $b_2$ , and  $b_3$  for Model 5**



**Figures 24 - 26: Posterior density plots of  $b_1$ ,  $b_2$ , and  $b_3$  for Model 6**



**Figures 27 - 29: Posterior density plots of  $b_1$ ,  $b_2$ , and  $b_3$  for Model 7**



**Figures 30 - 32: Posterior density plots of  $b_1$ ,  $b_2$ , and  $b_3$  for Model 8**

These density plots appear to be smooth and unimodal for most cases.



## DISCUSSION

This paper discusses a random changepoint model for HIV immunological markers, including CD4, CD8, and viral RNA levels at times surrounding the initiation of HAART therapy. The model was fitted in WinBUGS using a Bayesian approach, which allows flexibility in model specification. We chose noninformative priors for the model parameters and variance components. The model selection was directed by use of DIC, which combines a measure of goodness-of-fit as well as a measure of model complexity. DIC has limitations, but in our case provided credible values.

Our analysis confirms previous findings that in subjects with viral suppression the CD4, CD8 and viral RNA loads show a two-step change in their levels (an increase for CD4 and CD8 cell counts, and a decrease for viral RNA). We can improve prediction by allowing the changepoint of the rebound to vary between subjects. This also reinforces the idea that the T-cell rebound process may follow the viral load process, where the steep decline lasts for two weeks after beginning treatment (Ghosh and Vaida, 2006). Future research into the joint modeling of T-cell counts and RNA loads may be needed to further describe the association between the two.

Abrams, Donald I., Goldman, Anne I., Launer, Cynthia, Korvick, Joyce A., Neaton, James D., Crane, Lawrence R., Grodesky, Michael, Wakefield, Steven, Muth, Katherine, Kornegay, Sandra, Cohn, David L., Harris, Allen, Luskin-Hawk, Roberta, Markowitz, Norman, Sampson, James H., Thompson, Melanie, Deyton, Lawrence, and The Terry Bein Community Programs for Clinical Research on AIDS (1994). "A comparative trial of didanosine or zalcitabine after treatment with zidovudine in patients with human immunodeficiency virus infection." New England Journal of Medicine **330**: 657-662.

Berg, Andreas, Meyer, Renate, and Yu, Jun (2004). "Deviance Information Criterion for Comparing Stochastic Volatility Models." Journal of Business & Economic Statistics **22**(1): 107-121.

Deeks, S. G., Kitchen, C.M.R., Liu, Lea, Guo, Hua, Gascon, Ron, Narváez, Amy B., Hunt, P., Martin, J.N., Kahn, J.O., Levy, Jay, McGrath, Michael S., and Hecht, Frederick M. (2004). "Immune activation set point during early FHV infection predicts subsequent CD4(+) T-cell changes independent of viral load." Blood **104**(4): 942-947.

Ghosh, P. and F. Vaida (2006). "Random Changepoint Modeling of HIV Immunologic Responses."

Guo, X. and Carlin, B.P. (2004), "Separate and joint modeling of longitudinal and event time data using standard computer packages." The American Statistician, **58**:16-24.

Gumel, A. A. (2001). "A mathematical model for the dynamics of HIV-1 during the typical course of infection." Nonlinear Analysis, **47**(3): 1773-1783.

Halpern, Aaron (2000). "Multiple changepoint testing for an alternating segments model of a binary sequence." Biometrics, **56**: 903-908.

Hunt, Peter W., Deeks, Steven G., Rodriguez, Benigno, Valdez, H., Shade, Starley B., Abrams, D.I., Kitahata, Mari M., Krone, M., Neilands, T.B., Brand, R.J., Lederman, M.M., and Martin, J.N. (2003). "Continued CD4 cell count increases in HIV-infected adults experiencing 4 years of viral suppression on antiretroviral therapy." AIDS, **17**(13): 1907-1915.

Kimball, John W. (2005). "Retroviruses." Available at <http://users.rcn.com/jkimball.ma.ultranet/BiologyPages/R/Retroviruses.html>.

Kiuchi, Amy S., Hartigan, J.A., Halford, Theodore R., Rubenstein, Pablo, and Stevens, Cladd E. (1995). "Change points in the series of T4 counts prior to AIDS." Biometrics, **51**: 236 - 248.

Lange, Nicholas, Carlin, Bradley P., and Gelfand, Alan E. (1992). "Hierarchical Bayes models for the progression of HIV infection using longitudinal CD4 T-cell numbers." Journal of the American Statistical Association, **87**(419): 615-626.

Loschi, R.H., Cruz, F.R.B., and Arellano-Valle, R.B. (2005). "Multiple change point analysis for the regular exponential family using the product partition model." Journal of Data Science **3**: 305-330.

Myung, Jay I., Karabatsos, George, and Iverson, Geoffrey J. (2005). "A Bayesian approach to testing decision making axioms." Journal of Mathematical Psychology **49**: 205-225.

Taylor, Wayne A (2000). "Change-Point analysis: A powerful new tool for detecting changes." Available at <http://www.variation.com/cpa/tech/changepoint.html>.

Wang, L. C. (2006). "Mathematical analysis of the global dynamics of a model for HIV infection of CD4(+) T cells." Mathematical Biosciences **200**(1): 44-57.

Wikipedia Contributors. "HIV." Available at <http://en.wikipedia.org/wiki/Hiv>.

Wikipedia Contributors. "T Cell." Available at [http://en.wikipedia.org/wiki/T\\_cells](http://en.wikipedia.org/wiki/T_cells).

## APPENDIX A

### Implementation Using WinBUGS

Following is a copy of the WinBUGS code to implement the methods described in this paper.

*When  $K_i$  follow the exponential distribution:*

```
model{

# Level 1: Longitudinal data with changepoint K[i]

for (i in 1:N) {
  for (j in 1:M) {
    Y[i,j] ~ dnorm(muy[i,j], tauY)
    muy[i,j] <- beta[1]+ beta[2]*(t[j]-K[i])*min(K[i],t[j]) + beta[3]*(t[j]-K[i])*step(t[j]-K[i]) + b[i]

  }

}

# Level 2: random effects
b[i] ~ dnorm(0,invSigma)

K[i] ~ dexp(lambda)

}

# Level 3: priors
beta[1:3]~dmnorm(zero3[],invSigmabeta[,])

tauY~dgamma(0.1,0.1)
#tauYnew~dgamma(0.1,0.1)
invSigma ~ dgamma(0.1,0.1)

# derived variables
for (i in 1:N){ for(j in 1:M){
  res[i,j] <- Y[i,j]-muy[i,j]
  tt[i,j] <- t[j]
}}
sigmaY <- 1/sqrt(tauY)
sigmab <- 1/sqrt(invSigma)
meanK <- mean(K[])
}
```

*When  $K_i$  follow the truncated normal distribution:*

```

model{

# Level 1: Longitudinal data with changepoint K[i]

for (i in 1:N) {
  for (j in 1:M) {
    Y[i,j] ~ dnorm(muy[i,j], tauY)
    muy[i,j] <- beta[1]+ beta[2]*(t[j]-K[i])*min(K[i],t[j]) + beta[3]*(t[j]-K[i])*step(t[j]-K[i]) + b[i]

  }

}

for ( i in 1:N){
# Level 2: random effects
b[i] ~ dnorm(0,invSigma)

K[i] ~ dnorm(2,0.1)|(0, )

}

# Level 3: priors
beta[1:3]~dmnorm(zero3[],invSigmabeta[,])

tauY~dgamma(0.1,0.1)
#tauYnew~dgamma(0.1,0.1)
invSigma ~ dgamma(0.1,0.1)

# derived variables
for (i in 1:N){ for(j in 1:M){
res[i,j] <- Y[i,j]-muy[i,j]
tt[i,j] <- t[j]
}}
sigmaY <- 1/sqrt(tauY)
sigmab <- 1/sqrt(invSigma)
meanK <- mean(K[])
}

```